



REC'D 18 OCT 2004
WIPO PCT

PRIORITY DOCUMENT

SUBMITTED OR TRANSMITTED IN
COMPLIANCE WITH RULE 17.1(a) OR (b)

Patent Office
Canberra

I, JULIE BILLINGSLEY, TEAM LEADER EXAMINATION SUPPORT AND SALES hereby certify that annexed is a true copy of the Provisional specification in connection with Application No. 2003905362 for a patent by PROTEOME SYSTEMS INTELLECTUAL PROPERTY PTY LTD as filed on 01 October 2003.



WITNESS my hand this
Twelfth day of October 2004

JULIE BILLINGSLEY
TEAM LEADER EXAMINATION
SUPPORT AND SALES

BEST AVAILABLE COPY

AUSTRALIA

Patents Act 1990

Proteome Systems Intellectual Property Pty Ltd

PROVISIONAL SPECIFICATION

Invention Title:

A method for determining the biological likelihood of theoretical compositions or structures

The invention is described in the following statement:

Field of the Invention

This invention relates to a method of determining the biological likelihood of theoretical compositions or structures, particularly glycans.

5 Background of the Invention

Glycans (sugar structures/oligosaccharides) are usually composed of varying numbers of less than a dozen biologically-occurring monosaccharides. When considered purely in terms of their masses there are usually only about 3-6 different mass-unique monosaccharides in a typical glycan structure. The most frequently 10 encountered unique-mass monosaccharides are Hex (mass 162 Da; includes all hexose monosaccharides), HexNAc (mass 203 Da; includes all acetamidohexose monosaccharides), dHex (mass 146 Da; includes all deoxyhexose monosaccharides), Pent (mass 132 Da; includes all pentose monosaccharides), and NeuAc (mass 291 Da; N-acetylneuraminic (sialic) acid). There are several other biologically extant, though 15 less-frequently encountered component monosaccharides, such as KDN, HexA, NeuGc. Other non-monosaccharide adducts such as sulfate (S; mass 79.97 Da), phosphate (P; mass 97.98 Da), methyl (14 Da), and acetyl (42 Da) are also occasionally observed on biologically-occurring oligosaccharides.

It is often the case during the characterisation of biological molecules that a 20 precise mass may be ascertained for each biological molecule but its composition and identity are unknown. Given a reasonably accurate mass, such as would normally be obtained by mass spectrometry, the monosaccharide composition of an unknown glycan can be theorised by determining, by computation, the set of monosaccharide compositions that are within a reasonable mass deviation (or tolerance) of the observed 25 mass. This approach forms the basis of glycomod (<http://us.expasy.org/tools/glycomod/>) a publicly available research tool. The shortcomings of this tool, and of this purely theoretical approach by extension, is that a large number of compositions are returned for any mass of larger than moderate size, and that the majority (90-99%) of these have little in common with known biologically 30 extant compositions.

The aim of the present invention is to attempt to alleviate some of the above described problems and to reduce the large number of irrelevant compositions returned by existing tools.

Any discussion of documents, acts, materials, devices, articles or the like which 35 has been included in the present specification is solely for the purpose of providing a context for the present invention. It is not to be taken as an admission that any or all of

these matters form part of the prior art base or were common general knowledge in the field relevant to the present invention as it existed before the priority date of each claim of this application.

5 **Summary of the Invention**

In a first broad aspect, the present invention incorporates statistical measures of biological relevance for the theoretical compositions returned.

Typically, biological relevance, expressed as a numerical score, or biological index, is determined by statistical comparison to an established reference set of known 10 and fully characterised compositions, in the case of glycans a reference set such as the Glycosuite (<http://www.glycosuite.com>) database. The biological index of any given composition may then be used as a basis for discarding biologically "unlikely" compositions, as well as for ranking (sorting) of returned compositions by biological likeliness.

15 Empirically, for glycans this allows between 90-99.9% of theoretical compositions returned by any given search to be discarded, whilst preserving and ranking the remaining, biologically likely compositions.

In one aspect the present invention provides a method of determining the likelihood of a theoretical candidate composition comprising the steps of:

20 selecting a reference group of known characterised compositions;

establishing statistical characteristics relating to components of or other features of the known characterised composition;

comparing the statistical characteristics of the known characterised compositions with corresponding components or features in the theoretical candidate 25 compositions to establish a likelihood of those compositions occurring.

More particularly, for glycans, in one aspect the present invention provides a method of characterising glycans comprising the steps of:

providing a search mass of a glycan whose composition is to be determined;

generating a list of theoretical glycans made up of components, including 30 monosaccharides, whose total mass is within a predetermined tolerance of the search mass;

selecting a reference group of known characterised glycan compositions of approximately similar mass to the search mass;

establishing the mean and standard deviation of each component appearing in 35 the reference group of the known characterised glycan compositions;

for each theoretical glycan candidate calculating a partial score for each component in that theoretical glycan candidate based on the difference between the observed number of the component in the theoretical glycan candidate and the mean for that component in the reference group, divided by the standard deviation;

5 combining the partial scores to provide an indication of the likelihood of that theoretical glycan candidate occurring.

The partial scores may be combined in any suitable manner. One way, for example, is by multiplying the partial scores together.

By the use of actual biological information, the present invention is able to
10 discern biologically likely compositions from the vast majority of compositions of similar mass, but whose compositions differ greatly from known, biologically extant compositions. For example, for glycans where the publicly available web tool glycomod returns over 100 theoretical compositions for the mass 1300 Da +/- 0.5 Da, a tool embodying the present invention returns 2 biologically likely compositions and
15 109 biologically unlikely compositions (which would normally be discarded).

Although the main application of the present invention is to the delineation of biologically likely and unlikely sugar compositions for the purposes of sugar structure/composition elucidation, the generic methodology of using known biological data as a means to refine, interpret, and/or rank theoretical or empirical data may be
20 used for many other applications.

Detailed Description of a Preferred Embodiment

The present invention is implemented on a computer means running software carrying out the algorithms and process of the method.

The first input to a search using a method to determine a glycan composition
25 (the "search glycan") embodying the present invention, is a search mass (which is typically in Daltons). The search mass is typically an empirically determined mass of the "search glycan" which is to be characterised determined by mass spectrometry or other means i.e. the mass of the search glycan whose composition is to be determined.

A search mass tolerance (in Daltons) is also input. Typically this will be a
30 relatively small value depending on the expected accuracy of the empirically determined search mass and typically may be of the order of +/- 0.1Da. Also input is a "biological index" cut-off. The biological index is a measure of a theoretical glycan composition's likelihood and its derivation is explained in more detail below. The cut off is the value of that index above which candidate compositions are discarded as
35 being too unlikely to occur in the real world. Also input is a "maximum composition"

which indicates the maximum allowable number of each monosaccharide in each theoretical glycan composition. By way of example, if it were known for a fact that the glycan to be characterised contained no sialic acid, the theoretical glycan compositions generated as potential matches for the search mass would also exclude sialic acid. This 5 reduces the amount of computation required and improves speed and accuracy. In the system implementing the method, defaults would typically be provided for those inputs, except of course for the search mass.

Other optional parameters may also be exposed to the user to further modify the performance of the search. The output of the composition search is a list of candidate 10 theoretical glycan compositions, whose mass is within the search mass tolerance of the search mass, and whose biological index is less than the biological index cut-off. In theory one of those candidates matches the composition of the search glycan.

The composition search is performed as follows:

Reference statistics for the given search mass are determined from the 15 (Glycosuite) database. This process is described in more detail below.

Monosaccharides are recursively recombined in varying numbers such that every possible combination of allowed monosaccharides is created. Compositions whose mass does not fall within the search mass tolerance are discarded, as are compositions for which the number of any monosaccharide exceeds the maximum 20 number of that monosaccharide specified by the "maximum composition". The result is a list of theoretical candidate glycan compositions.

The biological index of candidate compositions is determined as described below. Compositions whose biological index does not satisfy the biological index cut-off are discarded.

25 The remaining compositions are presented to the user in order of biological index. Typically the list will be short and may only include one or two candidates. This compares with the hundreds of candidates typically produced by Glycomod, each of which has to be individually reviewed and assessed.

Calculation of Biological Index

30 Inputs to the process are a composition, and a reference data set of known sugar compositions/structures. The reference set may be from any suitable database or data source such as Glycosuite. The output of the process is a numerical biological index.

The determination of biological index for a given search glycan composition proceeds as follows:

35 The mass of the composition is the search mass or may be determined by the sum of the residue masses of each monosaccharide/component in the composition.

By reference to the reference set of known biological compositions, the mean and standard deviation of every monosaccharide/component in the database within an arbitrary mass range (eg: +/- 200 Da) of the mass of the composition is determined. Obtaining statistics from a range of masses around the given composition's mass is

5 necessary in order to obtain a sufficiently large sample size (preferably at least 100 known compositions). In the case of the Glycosuite database of known sugar structures, a mass tolerance of 200 Da was empirically determined to be sufficient to provide in excess of 100 known compositions for search masses up to around 3500.

By way of example if the search mass were 1000 Da there may be 100 known

10 glycans in the database whose mass is between 800 and 1200 Da. The mean and standard deviation of each of every monosaccharide/component appearing in those known glycans in the database is then determined. If we take HexNAc as an example we may find that, on average, the 100 known glycans contain 3.3 HexNAc monosaccharides with a standard deviation of 2.3. This process is repeated to calculate

15 the mean and standard deviation for each monosaccharide component Hex, dHex, pent et al, and each adduct in the known glycans, if adducts are being accounted for.

For each theoretical candidate glycan composition "Partial scores" are then determined from the means and standard deviations calculated above. These are calculated for each monosaccharide in the given composition as the absolute value of

20 the difference between the mean number of that monosaccharide in the reference set and the observed number of that monosaccharide in the theoretical candidate composition, divided by the standard deviation of that monosaccharide in compositions from the reference set. ie:

$$25 \quad partialscore_{monosac} = \frac{|mean_{monosac} - observed_{monosac}|}{stdev_{monosac}}$$

where $mean_{monosac}$ is the mean number of the given monosaccharide in the reference data set (Glycosuite); $mean_{monosac}$ is the number of the given monosaccharide in the theoretical candidate composition; and $stdev_{monosac}$ is the standard deviation of the given monosaccharide in the reference data set.

By way of example if the theoretical glycan composition includes two HexNAc, three Hex and 1 NeuAc, the partial score for each of those three monosaccharides is calculated for that theoretical candidate glycan composition. Partial scores need not be calculated for monosaccharides which do not appear in the candidate theoretical glycan composition.

In the event that the $mean_{monosac}$ equals the $mean_{monosac}$ for a particular glycan, the system is arranged to give the partial score a minimum value of 0.01.

Thus, the partial score of a monosaccharide is in fact the number of standard deviations the number of away from the mean that that monosaccharide is in the theoretical candidate composition. In a normal distribution, approximately 68% of all data points lie within 1 standard deviation of the mean, ~93% within 2 standard deviations, over 99% within 3. Assuming that the distributions of monosaccharide number for the mass range used to obtain the initial means and standard deviations for the given search mass are sufficiently close to normal, then partial scores of 3 or less for any monosaccharide indicate that the number of those monosaccharides are within 99% of all compositions of similar mass in Glycosuite.

- 5
- 10

Partial scores are then combined in some manner to derive a single numeric score; this being the biological index. The actual mathematical derivation of the biological index may be arrived at using multiple means; different formulae exhibit subtle differences in their sensitivity to large partial scores and other criteria. For this reason, biological index for the purposes of the present invention may be considered merely as a numerical value that is representative of, and derived from, the magnitudes of the differences between a given composition and a population of known compositions of a similar mass. Presently, a biological index is calculated from partial scores as the product of all the partial scores from the theoretical candidate composition; ie:

- 15
- 20

$$BI = \prod_{monosac_0}^{monosac_n} partialscore_{monosac}$$

The Biological Index is adept at excluding very poor matches but at the same time if a candidate theoretical glycan composition has a very large (i.e. poor) partial score for one monosaccharide but low partial scores for the other monosaccharide components, the candidate may have an acceptably low Biological Index hence the system does not discard candidates which have only one poor partial score.

- 25

The process of calculating the partial scores is carried out for each theoretical glycan composition as discussed above. Compositions whose biological index does not satisfy the biological index cut-off are discarded. The remaining compositions are presented to the user in order of biological index. Typically the list will be short and may only include one or two candidates. This compares with the hundreds of candidates typically produced by Glycomod, each of which has to be individually reviewed and assessed.

- 30

The key element of the present invention is the use of biological data as a means to score the quality of theoretical data, in this case, sugar compositions. The actual manner in which a biological score/index is calculated is largely arbitrary; different formulae for calculating a biological index exhibit different characteristics with respect

5 to their tolerance to large compositional differences from the determined mean, and in their propensity to extrapolate the compositions present in the reference database.

Although the present invention as described above is concerned with the use of known sugar structures/compositions as a means to discern/elucidate monosaccharide composition given only a mass, the concept could be extended to other compositions

10 and to the use of other structural characteristics, for example linkage and branching, as reference data for determining and/or ascertaining the quality of complete sugar structures for other investigative techniques, such as glycan fragment mass fingerprinting (see the applicant's co-pending provisional patent application No 2003902907, the entire contents of which are incorporated herein by reference).

15 It will be appreciated by persons skilled in the art that numerous variations and/or modifications may be made to the invention as shown in the specific embodiments without departing from the spirit or scope of the invention as broadly described. The present embodiments are, therefore, to be considered in all respects as illustrative and not restrictive.

Dated this first day of October 2003

Proteome Systems Intellectual Property Pty
Ltd
Patent Attorneys for the Applicant:

F B RICE & CO